

Lab 5: MNIST/FMNIST

2024-02-29

Homework (lab 4 or lab 5)

The homework consists of questions to be solved using R or python commands. You can choose between lab 4 and lab 5. You will type the necessary commands to solve the questions and the necessary comments to understand your calculations in a pdf/rtx file that you will save under the name TML_yourname. Your homework is to be sent via EPI before Monday February 26, 5pm (before class).

The MNIST database is a well-known database of handwritten digit images (0-9) that is used to evaluate machine learning algorithms. The database contains 60,000 training images and 10,000 test images. The MNIST data is stored in 4 binary files, which can be difficult to manage directly. Download the data (directly on Yann LeCun's website [here](#) or [here](#) for Fashion-MNIST) and you can use the code below to load them under R. In python, they (MNIST and FMNIST) are accessible via the tensorflow API for example.

```
rm(list = ls())

library(foreach)
library(doparallel)

## Ladowanie wymaganego pakietu: iterators

## Ladowanie wymaganego pakietu: parallel

library(e1071)
library(e1071)
library(ggplot2)

setwd("C:/Machine_Learning_A/Final_Project")

nTrain<- readBin(con = "train-images.idx3-ubyte", what = "integer", n = 4,
                size = 4, endian = "big")[2]
nTest<- readBin(con = "t10k-images.idx3-ubyte", what = "integer", n = 4,
               size = 4, endian = "big")[2]
nRows<- readBin(con = "train-images.idx3-ubyte", what = "integer", n = 4,
               size = 4, endian = "big")[3]
nCols<- readBin(con = "train-images.idx3-ubyte", what = "integer", n = 4,
               size = 4, endian = "big")[4]

# Load training data
images.train<- array(as.numeric(readBin(con = "train-images.idx3-ubyte",
                                       what = "raw",
                                       n = nTrain * nRows * nCols + 16,
                                       endian = "big"))[-(1:16)]),
                    dim = c(nCols, nRows, nTrain))
labels.train<- as.numeric(readBin(con = "train-labels.idx1-ubyte",
                                 what = "raw",
                                 n = nTrain + 0,
                                 endian = "big"))[-(1:9)]

# Load test data
images.test<- array(as.numeric(readBin(con = "t10k-images.idx3-ubyte",
                                       what = "raw",
                                       n = nTest * nRows * nCols + 16,
                                       endian = "big"))[-(1:16)]),
                    dim = c(nCols, nRows, nTest))
labels.test<- as.numeric(readBin(con = "t10k-labels.idx1-ubyte",
                                 what = "raw",
                                 n = nTest + 8,
                                 endian = "big"))[-(1:9)]
```

0. Definitions of function used later on:

```
set.seed(12)

#Functions dataPrep and createMatrix are used to manipulate data (sample etc.)
dataPrep<- function(p, n, Data_Images, Data_Labels){ #p is portion of Data to keep, n is the pa
#Labels to keep

sample_size<- length(Data_Images[1,1]) * p

small_sample<- sample(length(Data_Images[1, 1]),
                     size = sample_size)

Images<- Data_Images[,small_sample]

Labels<- Data_Labels[small_sample]

subs_train<- which(Labels<= n) #Choosing data which has label 0 1 or 2
Images<- Images[,subs_train]
Labels<- Labels[subs_train]

ret_list<- list(Images, Labels)

return(ret_list) #Return is an array with first element with Images and the second with Labels
}

createMatrix<- function(Data_List){
MatrixTrain<- matrix(0, nrow = length(Data_List[[1]][1,1]),
                    ncol = length(c(Data_List[[1]][1,1])),

for(j in 1:length(Data_List[[1]][1, 1])){
for(i in 1:ncol(Data_List[[1]][1, 1])){
MatrixTrain[i,j]<- c(Data_List[[1]][i, , j])
}
}

MatrixTrain<- MatrixTrain/255

Labels_Train<- as.factor(Data_List[[2]])

ret_list<- list(Matrix_Train, Labels_Train)

return(ret_list)
}

#Functions fitAndTest and benchmarkValues are used to compare different parameters of a models.

fitAndTest<- function(train_data, test_data, cost_in = 1, gamma_in = 0.05){

svm_fit<- e1071::svm(train_data[[1]],
                    train_data[[2]],
                    kernel = "radial",
                    cost = cost_in,
                    gamma = gamma_in,
                    scale = TRUE)

ypred<- predict(svm_fit, test_data[[1]])

conf_mat<- table(predict = ypred, truth = test_data[[2]])

corr_pred<- 0

incorr_pred<- 0

for(i in 1:ncol(conf_mat)){
for(j in 1:ncol(conf_mat)){
if(i == j){
corr_pred<- corr_pred + conf_mat[i,j]
}
else {
incorr_pred<- incorr_pred + conf_mat[i,j]
}
}
}

Total_MissClassRate<- corr_pred/(corr_pred + incorr_pred)

return(Total_MissClassRate)
}

benchmarkValues<- function(Data_Img_Train, Data_Lab_Train, Data_Img_Test, Data_Lab_Test, n){
#n is the number of folds we are considering
#(e.g. 0, 1, 2 -> n=2)

MissClassVector<- vector(mode = "double", length = n)

sample_train_inc<- ((1/2)-(1/30))/8

sample_test_inc<- ((1/2)-(1/10))/8

for(i in 2:n){

p_train<- 1/30 + sample_train_inc * (i-2)

Train_Prep<- dataPrep(p_train, 1, Data_Img_Train, Data_Lab_Train)

Train_Data<- createMatrix(Train_Prep)

p_test<- 1/10 + sample_test_inc * (i-2)

Test_Prep<- dataPrep(p_test, 1, Data_Img_Test, Data_Lab_Test)

Test_Data<- createMatrix(Test_Prep)

MissClassVector[i]<- fitAndTest(Train_Data, Test_Data)

}

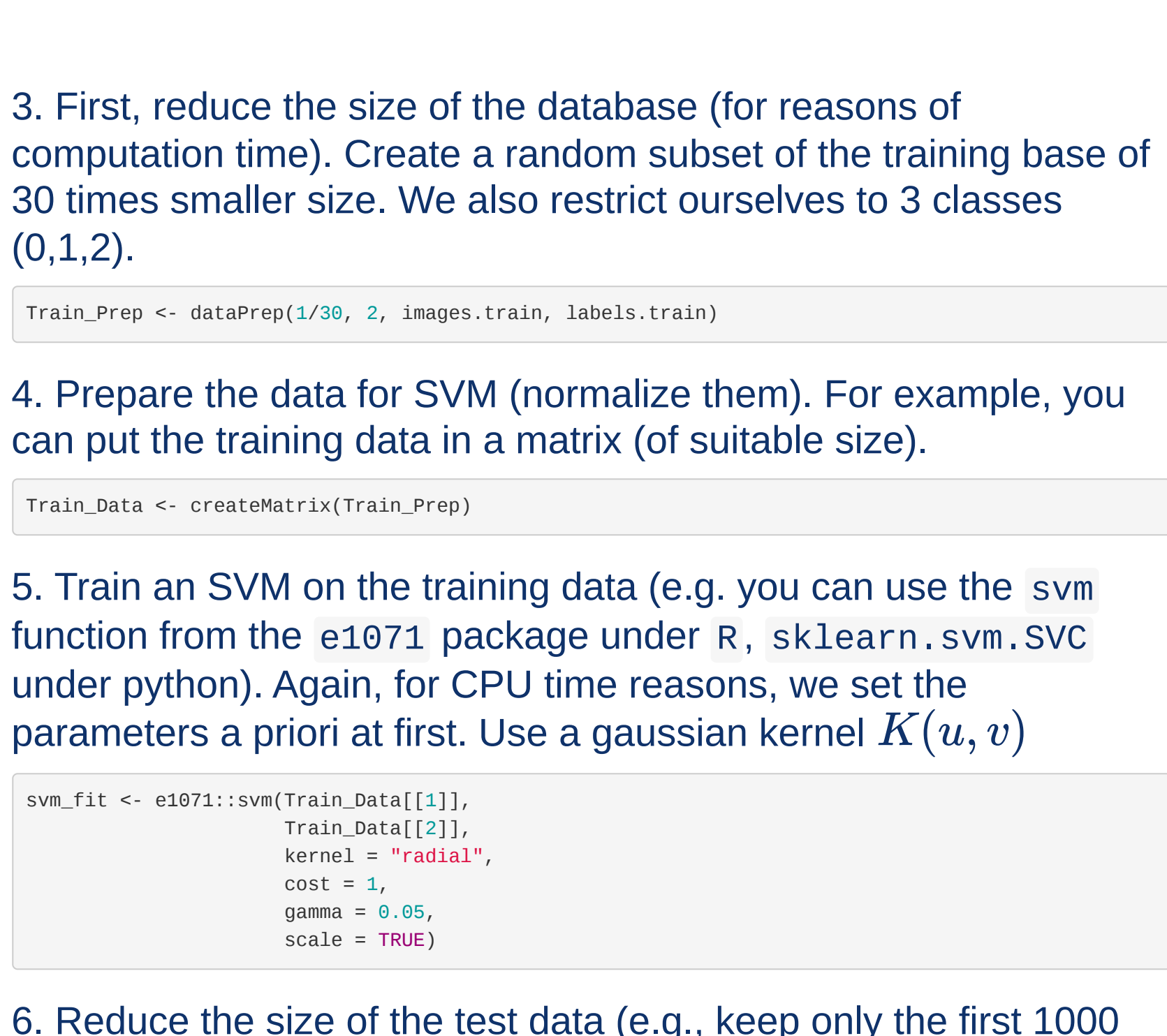
return(MissClassVector)
}
```

1. Display the first 10 images of the training base. You can use the functions `image` and `grey` or the package `imageR` (by converting them as `img` objects). With Python `matplotlib`.

```
par(mfrow = c(3, 3))

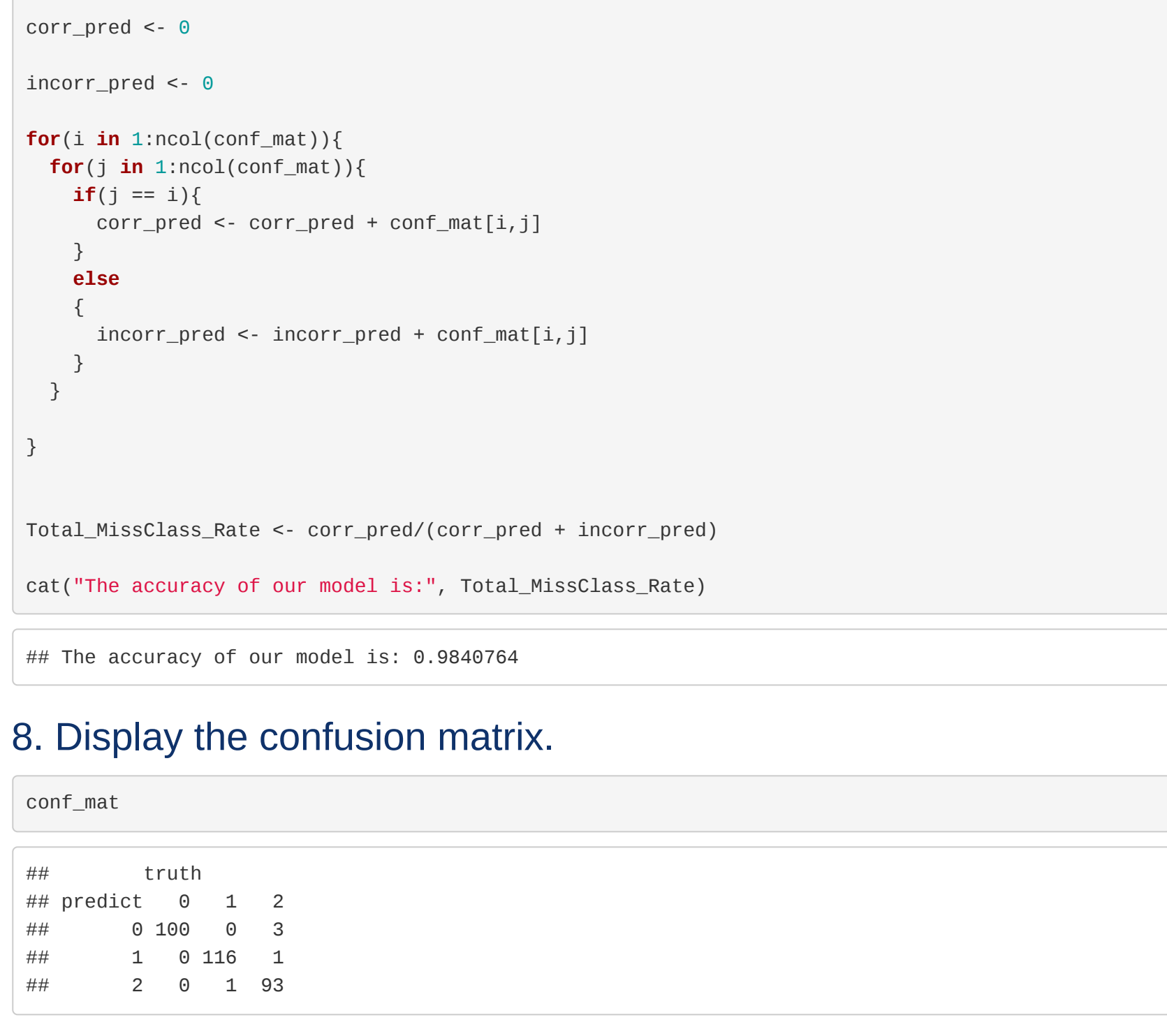
v<- c(1:28) # We use v to inverse the order of the rows (so that images are familiar to humans)

for(i in 1:9){
image(images.train[, i][, rev(v)], col = grey.colors(n = 255, rev = TRUE))
}
```



2. Same question for the first 10 images of the test base

```
par(mfrow = c(3, 3))
for(i in 1:9){
image(images.test[,i][, rev(v)], col = grey.colors(n = 255, rev = TRUE))
}
```



3. First, reduce the size of the database (for reasons of computation time). Create a random subset of the training base of 30 times smaller size. We also restrict ourselves to 3 classes (0,1,2).

```
Train_Prep<- dataPrep(1/30, 2, images.train, labels.train)
```

4. Prepare the data for SVM (normalize them). For example, you can put the training data in a matrix (of suitable size).

```
Train_Data<- createMatrix(Train_Prep)
```

5. Train an SVM on the training data (e.g. you can use the `svm` function from the `e1071` package under R, `sklearn.svm.SVC` under python). Again, for CPU time reasons, we set the parameters a priori at first. Use a gaussian kernel $K(u, v)$

```
svm_fit<- e1071::svm(Train_Data[[1]],
                    Train_Data[[2]],
                    kernel = "radial",
                    cost = 1,
                    gamma = 0.05,
                    scale = TRUE)
```

6. Reduce the size of the test data (e.g., keep only the first 1000 images) and prepare them in the same way as the training data.

```
Test_Prep<- dataPrep(0.1, 2, images.test, labels.test)
Test_Data<- createMatrix(Test_Prep)
```

7. Calculate the prediction error of the svm (on the first three classes only).

```
ypred<- predict(svm_fit, Test_Data[[1]])

conf_mat<- table(predict = ypred, truth = Test_Data[[2]])

corr_pred<- 0

incorr_pred<- 0

for(i in 1:ncol(conf_mat)){
for(j in 1:ncol(conf_mat)){
if(i == j){
corr_pred<- corr_pred + conf_mat[i,j]
}
else {
incorr_pred<- incorr_pred + conf_mat[i,j]
}
}
}

Total_MissClassRate<- corr_pred/(corr_pred + incorr_pred)

cat("The accuracy of our model is:", Total_MissClassRate)
```

The accuracy of our model is: 0.9840764

8. Display the confusion matrix.

```
conf_mat

##      truth
## predict 0 1 2
##      0 100 0 3
##      1 0 116 1
##      2 0 1 83
```

9. Compare the results using the classification method of your choice (e.g. CNN).

```
library(teras)

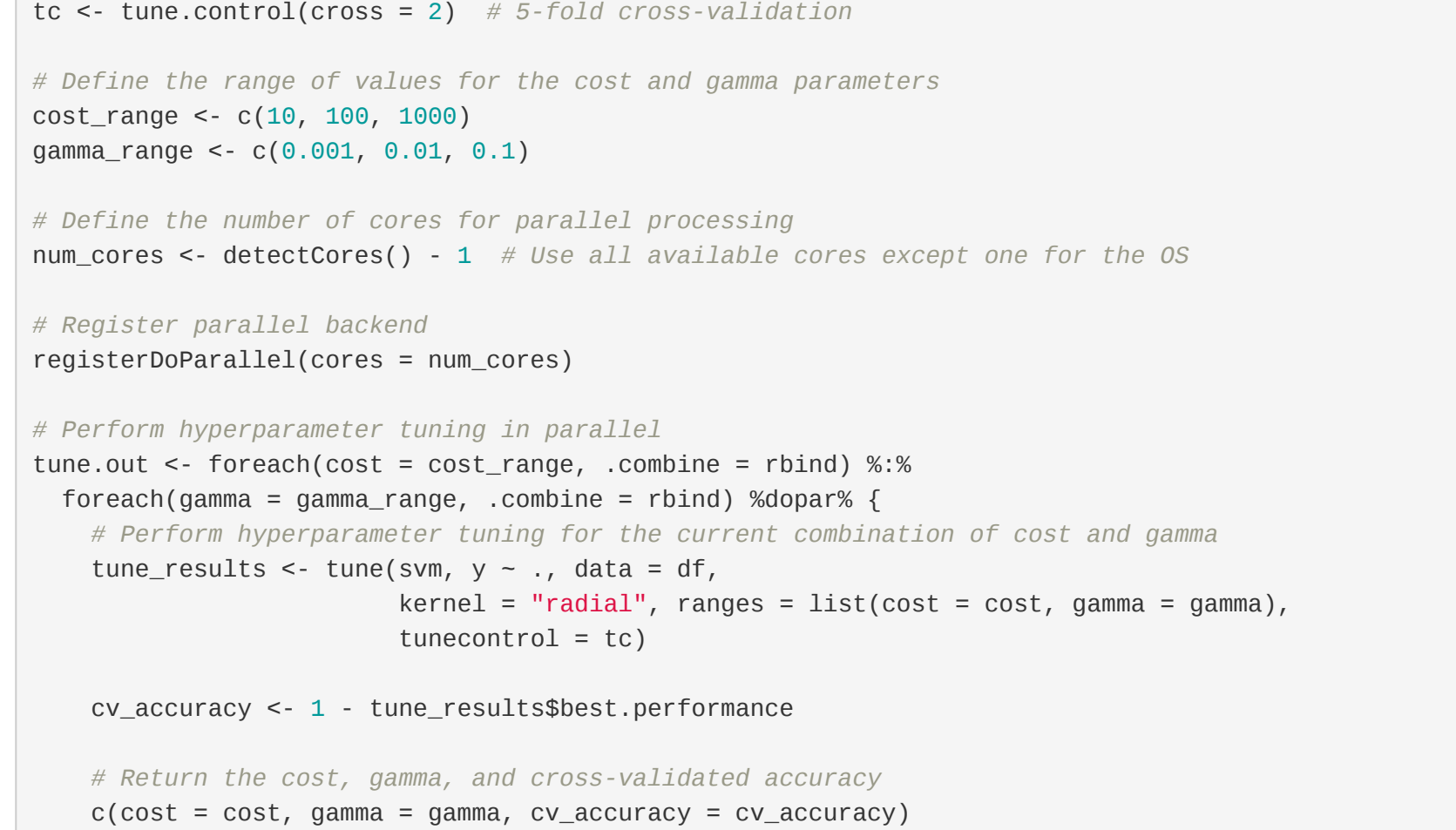
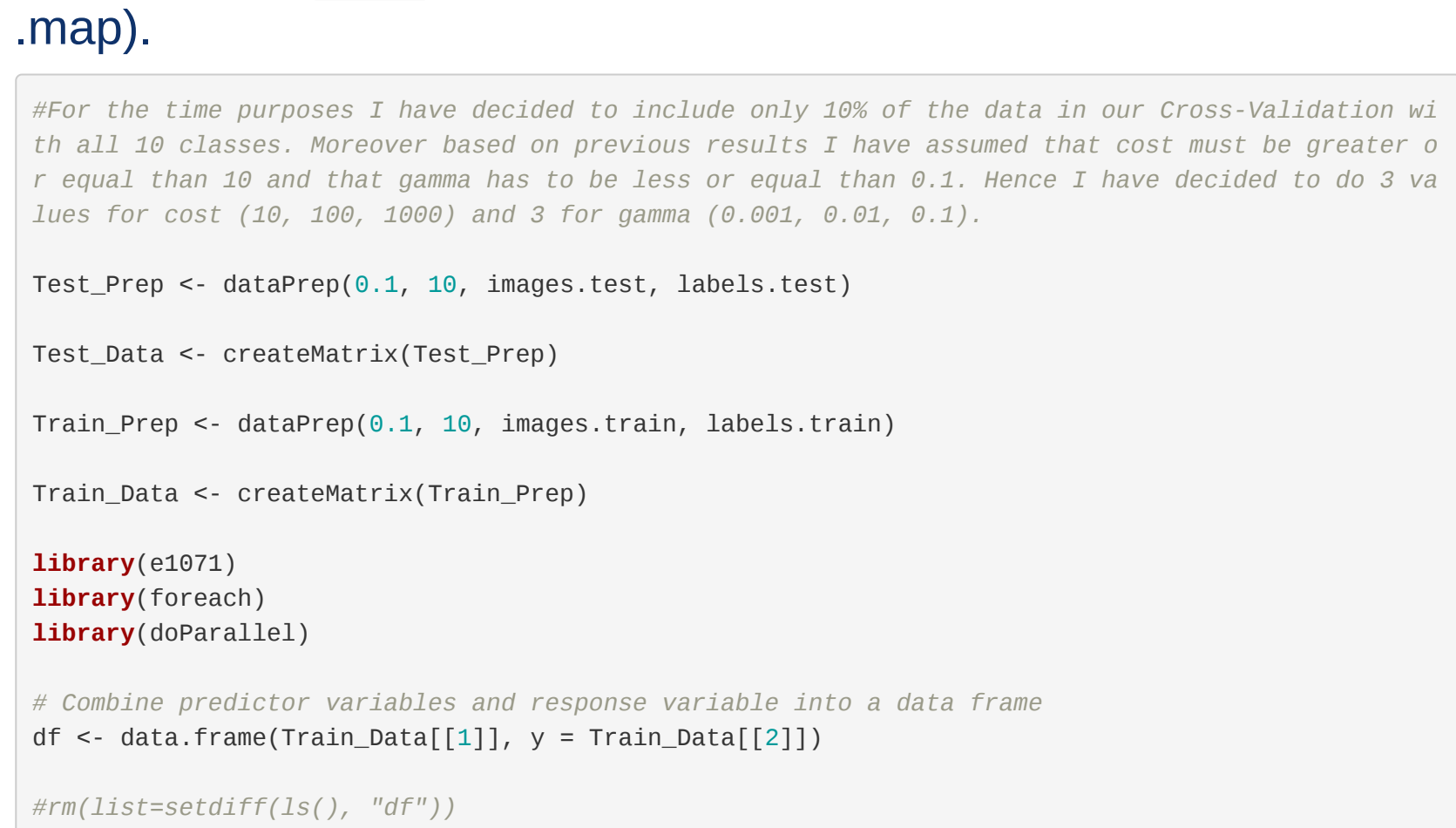
# Load the MNIST dataset
mist<- dataset.mnist()
x_train<- mist$trainx
y_train<- mist$trainy
x_test<- mist$testx
y_test<- mist$testy

# Preprocess the data
x_train<- array_reshape(x_train, c(nrow(x_train), 28, 28, 1))
x_test<- array_reshape(x_test, c(nrow(x_test), 28, 28, 1))
x_train<- x_train / 255
x_test<- x_test / 255
y_train<- to_categorical(y_train, 10)
y_test<- to_categorical(y_test, 10)

# Define the CNN model architecture
model<- keras_model_sequential() %>%
  layer_conv2d(filters = 32, kernel_size = c(3, 3), activation = 'relu',
               input_shape = c(28, 28, 1)) %>%
  layer_max_pooling2d(pool_size = c(2, 2)) %>%
  layer_conv2d(filters = 64, kernel_size = c(3, 3), activation = 'relu') %>%
  layer_max_pooling2d(pool_size = c(2, 2)) %>%
  layer_flatten() %>%
  layer_dense(units = 128, activation = 'relu') %>%
  layer_dropout(rate = 0.5) %>%
  layer_dense(units = 10, activation = 'softmax')

# Compile the model
model %>% compile(
  loss = 'categorical_crossentropy',
  optimizer = optimizer_adam(),
  metrics = c('accuracy')
)

# Train the model
history<- model %>% fit(
  x_train, y_train,
  epochs = 5,
  batch_size = 128,
  validation_split = 0.2
)
```



11. Using the packages `foreach` and `doMC` (you can consult the related vignettes `foreach` or `doMC`) parallelize the calibration step by K-fold cross-validation of the parameters γ and cost of the SVM. We will take for example $\gamma = \text{cost} = 10^{-3.3}$ and $K = 5$. In python, the `multiprocessing` module (<https://docs.python.org/2/library/multiprocessing.html>), and in particular the `Poo1` class (with its associated methods, `.apply`, `.map`).

```
#For the time purposes I have decided to include only 10% of the data in our Cross-validation wi
th all 10 classes. Moreover based on previous results I have assumed that cost must be greater o
r equal than 10 and that gamma has to be less or equal than 0.1. Hence I have decided to do 3 va
lues for cost (10, 100, 1000) and 3 for gamma (0.001, 0.01, 0.1).

Test_Prep<- dataPrep(0.1, 10, images.test, labels.test)

Test_Data<- createMatrix(Test_Prep)

Train_Prep<- dataPrep(0.1, 10, images.train, labels.train)

Train_Data<- createMatrix(Train_Prep)

library(e1071)
library(foreach)
library(doparallel)

# Combine predictor variables and response variable into a data frame
df<- data.frame(Train_Data[[1]], y = Train_Data[[2]])

#n=length(unique(df[, "y"]))

# Define the tuning control (e.g., cross-validation)
tc<- tune.control(cross = 2) # 5-fold cross-validation

# Define the range of values for the cost and gamma parameters
cost_range<- c(10, 100, 1000)
gamma_range<- c(0.001, 0.01, 0.1)

# Define the number of cores for parallel processing
num_cores<- detectCores() - 1 # Use all available cores except one for the OS

# Register parallel backend
registerDoParallel(cores = num_cores)

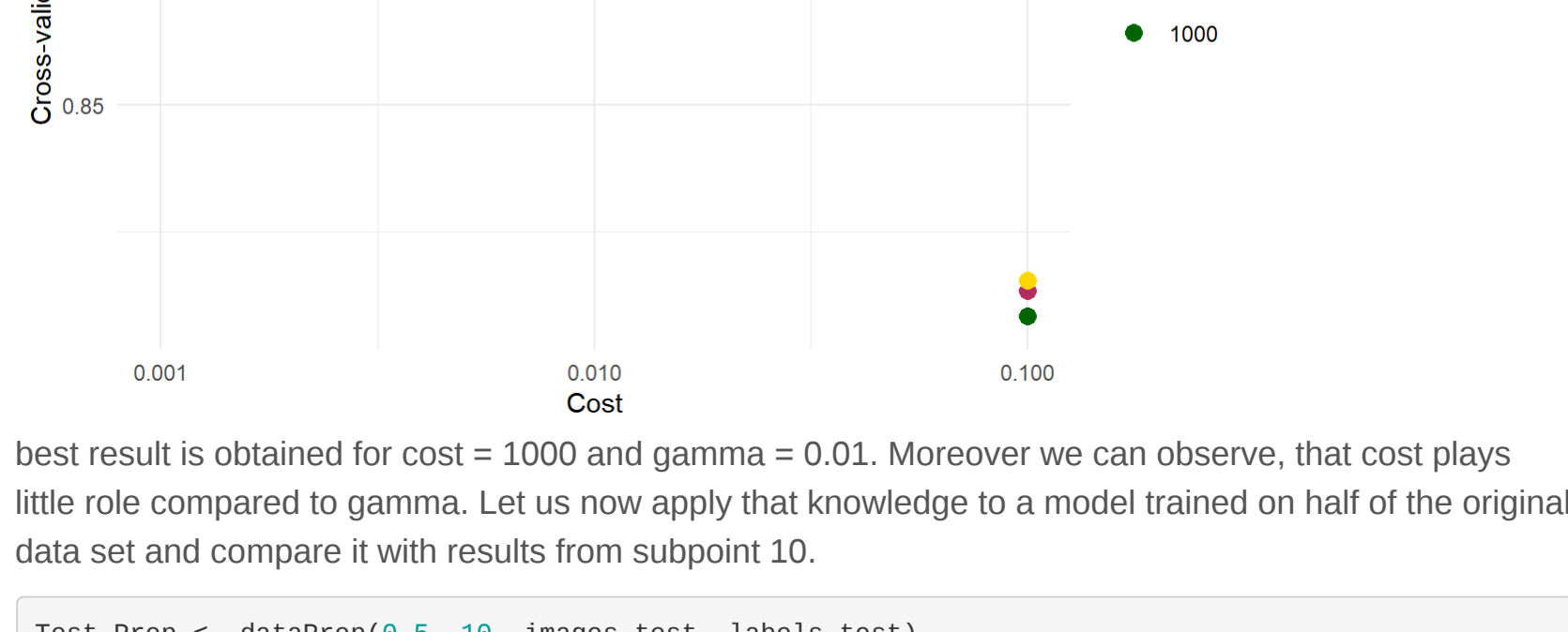
# Perform hyperparameter tuning in parallel
tune.out<- foreach(cost = cost_range, combine = rbind) %>%
  foreach(gamma = gamma_range, combine = rbind) %doParallel {
    # Perform hyperparameter tuning for the current combination of cost and gamma
    tune_results<- tune(svm, y ~ ., data = df,
                       kernel = "radial", ranges = list(cost = cost, gamma = gamma),
                       tunecontrol = tc)

    cv_accuracy<- 1 - tune_results$best.performance

    # Return the cost, gamma, and cross-validated accuracy
    c(cost = cost, gamma = gamma, cv_accuracy = cv_accuracy)
  }

# Stop parallel backend
stopImplicitCluster()

# Print the results
print(tune.out)
```



As we can see, the results for the model indicated by cross validation is doing slightly better than the old model. We may explain it as a very small difference in gamma, which seems to be the leading parameter here (as shown on the figures).